

A Bayesian nonparametric method for the LR assessment in case of rare type match

Giulia Cereda

October 8, 2015

Ingredients:

Ingredients:

- Crime case

Ingredients:

- Crime case
- Evidence (E)

Ingredients:

- Crime case
- Evidence (E)
- 2 Hypotheses of Interest: H_p vs H_d

Ingredients:

- Crime case
- Evidence (E)
- 2 Hypotheses of Interest: H_p vs H_d
- Background (B)

Ingredients:

- Crime case
- Evidence (E)
- 2 Hypotheses of Interest: H_p vs H_d
- Background (B)

$D = (E, B)$.

Ingredients:

- Crime case
- Evidence (E)
- 2 Hypotheses of Interest: H_p vs H_d
- Background (B)

$D = (E, B)$.

The court asks for the likelihood ratio

Ingredients:

- Crime case
- Evidence (E)
- 2 Hypotheses of Interest: H_p vs H_d
- Background (B)

$D = (E, B)$.

The court asks for the likelihood ratio

$$\frac{\Pr(H_p | D)}{\Pr(H_d | D)} = \underbrace{\frac{\Pr(D | H_p)}{\Pr(D | H_d)}}_{LR} \frac{\Pr(H_p)}{\Pr(H_d)}$$

Example

Ingredients:

Example

Ingredients:

- Crime case: murder

Ingredients:

- Crime case: murder
- Evidence (E): profile of the DNA trace found on the crime scene matches the suspect's DNA profile.

Ingredients:

- Crime case: murder
- Evidence (E): profile of the DNA trace found on the crime scene matches the suspect's DNA profile.
- 2 Hypotheses of Interest:

Ingredients:

- Crime case: murder
- Evidence (E): profile of the DNA trace found on the crime scene matches the suspect's DNA profile.
- 2 Hypotheses of Interest:
 - H_p : The suspect left the stain

Ingredients:

- Crime case: murder
- Evidence (E): profile of the DNA trace found on the crime scene matches the suspect's DNA profile.
- 2 Hypotheses of Interest:
 - H_p : The suspect left the stain
 - H_d : Someone else left the stain

Ingredients:

- Crime case: murder
- Evidence (E): profile of the DNA trace found on the crime scene matches the suspect's DNA profile.
- 2 Hypotheses of Interest:
 - H_p : The suspect left the stain
 - H_d : Someone else left the stain
- Background (B): database of DNA profiles from the population of possible perpetrators

DNA profiles

A DNA profile is a list of integers $h = (4 - 5 - 2 - 10)$ that code some characteristics in some portions of the DNA sequence of an individual: different persons can share the same profile.

A DNA profile is a list of integers $h = (4 - 5 - 2 - 10)$ that code some characteristics in some portions of the DNA sequence of an individual: different persons can share the same profile.

For H_p the match is a sure event,

A DNA profile is a list of integers $h = (4 - 5 - 2 - 10)$ that code some characteristics in some portions of the DNA sequence of an individual: different persons can share the same profile.

For H_p the match is a sure event,

For H_d the match is a random event with probability $p_h =$ frequency of the profile h of the suspect in the population of possible perpetrators.

DNA database

Database: a list of DNA profiles from a sample from the population of possible perpetrators

DNA database

Database: a list of DNA profiles from a sample from the population of possible perpetrators

DATABASE of size 10

Person 1 (4 – 10 – 6 – 7)

Person 2 (3 – 5 – 6 – 8)

Person 3 (3 – 7 – 8 – 10)

Person 4 (10 – 1 – 4 – 5)

Person 5 (3 – 7 – 8 – 10)

Person 6 (3 – 7 – 8 – 10)

Person 7 (1 – 5 – 7 – 2)

Person 8 (3 – 7 – 8 – 10)

Person 9 (3 – 5 – 6 – 8)

Person 10 (3 – 7 – 8 – 10)

DNA database

Database: a list of DNA profiles from a sample from the population of possible perpetrators

DATABASE of size 10

Person 1 (4 – 10 – 6 – 7)

Person 2 (3 – 5 – 6 – 8)

Person 3 (3 – 7 – 8 – 10)

Person 4 (10 – 1 – 4 – 5)

Person 5 (3 – 7 – 8 – 10)

Person 6 (3 – 7 – 8 – 10)

Person 7 (1 – 5 – 7 – 2)

Person 8 (3 – 7 – 8 – 10)

Person 9 (3 – 5 – 6 – 8)

Person 10 (3 – 7 – 8 – 10)

The database is used to find out the rarity of the matching profile.

LR assessment in the rare type match case

My research focuses on the LR assessment in the rare type match case, that is:

LR assessment in the rare type match case

My research focuses on the LR assessment in the rare type match case, that is:

- A match between the suspect's DNA profile and the crime stain's DNA profile.

LR assessment in the rare type match case

My research focuses on the LR assessment in the rare type match case, that is:

- A match between the suspect's DNA profile and the crime stain's DNA profile.
- This profile is not contained in the database B.

LR assessment in the rare type match case

My research focuses on the LR assessment in the rare type match case, that is:

- A match between the suspect's DNA profile and the crime stain's DNA profile.
- **This profile is not contained in the database B.**

Especially if the database is big, the profile seems to be rare.

LR assessment in the rare type match case

My research focuses on the LR assessment in the rare type match case, that is:

- A match between the suspect's DNA profile and the crime stain's DNA profile.
- This profile is not contained in the database B.

Especially if the database is big, the profile seems to be rare.

How rare?

Previous models

- Frequentist model:
(Cereda 2015) Frequentist approach to LR assessment in case of rare haplotype match
arXiv:1502.04083

- Frequentist model:
(Cereda 2015) Frequentist approach to LR assessment in case of rare haplotype match
arXiv:1502.04083
- Bayesian model:
(Cereda 2015) Full Bayesian approach to LR assessment in case of rare haplotype match
arXiv:1502.02406

- Frequentist model:
(Cereda 2015) Frequentist approach to LR assessment in case of rare haplotype match
arXiv:1502.04083
- Bayesian model:
(Cereda 2015) Full Bayesian approach to LR assessment in case of rare haplotype match
arXiv:1502.02406
- (Cereda 2015) Nonparametric Bayesian approach to LR assessment in case of rare haplotype match
arXiv:1506.08444

Assumption 1

There are so many different DNA types that they may be considered infinite.

Assumption 1

There are so many different DNA types that they may be considered infinite.

Parameter: $\mathbf{p} = (p_t | t \in T)$, T an infinite countable set, $p_t > 0$, $\sum p_t = 1$, to represent the (unknown) frequencies of all DNA types in Nature.

Assumption 1

There are so many different DNA types that they may be considered infinite.

Parameter: $\mathbf{p} = (p_t | t \in T)$, T an infinite countable set, $p_t > 0$, $\sum p_t = 1$, to represent the (unknown) frequencies of all DNA types in Nature.

Assumption 2

The particular list of integers that forms a DNA type is just a category: no structure assumed.

Assumption 1

There are so many different DNA types that they may be considered infinite.

Parameter: $\mathbf{p} = (p_t | t \in T)$, T an infinite countable set, $p_t > 0$, $\sum p_t = 1$, to represent the (unknown) frequencies of all DNA types in Nature.

Assumption 2

The particular list of integers that forms a DNA type is just a category: no structure assumed.

“DNA types” or “colors” is now the same.

Random partitions of $[n]$

Let $[n]$ denote the set $[n] = \{1, 2, \dots, n\}$.

Random partitions of $[n]$

Let $[n]$ denote the set $[n] = \{1, 2, \dots, n\}$.

A partition of the set $[n]$ will be denoted as $\pi_{[n]}$.

Random partitions of $[n]$

Let $[n]$ denote the set $[n] = \{1, 2, \dots, n\}$.

A partition of the set $[n]$ will be denoted as $\pi_{[n]}$.

Random partitions on the set $[n]$ will be denoted as $\Pi_{[n]}$.

DNA database can be reduced

DATABASE of size 10

Person 1 (4 – 10 – 6 – 7)

Person 2 (3 – 5 – 6 – 8)

Person 3 (3 – 7 – 8 – 10)

Person 4 (10 – 1 – 4 – 5)

Person 5 (3 – 7 – 8 – 10)

Person 6 (3 – 7 – 8 – 10)

Person 7 (1 – 5 – 7 – 2)

Person 8 (3 – 7 – 8 – 10)

Person 9 (3 – 5 – 6 – 8)

Person 10 (3 – 7 – 8 – 10)

DNA database can be reduced

DATABASE of size 10

Person 1 (4 – 10 – 6 – 7)

Person 2 (3 – 5 – 6 – 8)

Person 3 (3 – 7 – 8 – 10)

Person 4 (10 – 1 – 4 – 5)

Person 5 (3 – 7 – 8 – 10)

Person 6 (3 – 7 – 8 – 10)

Person 7 (1 – 5 – 7 – 2)

Person 8 (3 – 7 – 8 – 10)

Person 9 (3 – 5 – 6 – 8)

Person 10 (3 – 7 – 8 – 10)

DNA database can be reduced

DATABASE of size 10

Person 1 (4 – 10 – 6 – 7)

Person 2 (3 – 5 – 6 – 8)

Person 3 (3 – 7 – 8 – 10)

Person 4 (10 – 1 – 4 – 5)

Person 5 (3 – 7 – 8 – 10)

Person 6 (3 – 7 – 8 – 10)

Person 7 (1 – 5 – 7 – 2)

Person 8 (3 – 7 – 8 – 10)

Person 9 (3 – 5 – 6 – 8)

Person 10 (3 – 7 – 8 – 10)

Assumption 2 → data can be replaced by the equivalence classes on the indices of the relation “to have the same DNA type”.

This is a partition of the set $[n]$: $\{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}\}$

Data \mathcal{D} is made of the database + 2 new observations

Data \mathcal{D} is made of the database + 2 new observations

$\mathcal{D} = \pi_{[n+2]}$ partition of the set $\{1, 2, \dots, n + 2\}$

Data \mathcal{D} is made of the database + 2 new observations

$\mathcal{D} = \pi_{[n+2]}$ partition of the set $\{1, 2, \dots, n + 2\}$

Example:

Database $\rightarrow \pi_{[10]} = \{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}\}$

Data \mathcal{D} is made of the database + 2 new observations

$\mathcal{D} = \pi_{[n+2]}$ partition of the set $\{1, 2, \dots, n + 2\}$

Example:

Database $\rightarrow \pi_{[10]} = \{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}\}$

$\mathcal{D} \rightarrow \pi_{[12]} = \{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}, \{11, 12\}\}$

Data \mathcal{D} is made of the database + 2 new observations

$\mathcal{D} = \pi_{[n+2]}$ partition of the set $\{1, 2, \dots, n + 2\}$

Example:

Database $\rightarrow \pi_{[10]} = \{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}\}$

$\mathcal{D} \rightarrow \pi_{[12]} = \{\{1\}, \{2, 9\}, \{3, 5, 6, 8, 10\}, \{4\}, \{7\}, \{11, 12\}\}$

We can see the data as a random variable. In that case,

$$\mathcal{D} = \Pi_{[n+2]}.$$

The distribution of $\mathcal{D} = \Pi_{[n+2]}$ depends on \mathbf{p} . However, it does not depend on the order of the p_j .

The distribution of $\mathcal{D} = \Pi_{[n+2]}$ depends on \mathbf{p} . However, it does not depend on the order of the p_j .



We can consider directly the ordered vector

$$\mathbf{p} \in \nabla_\infty = \{(p_1, p_2, \dots), p_1 \geq p_2 \geq \dots > 0, \sum p_i = 1\}.$$

The distribution of $\mathcal{D} = \Pi_{[n+2]}$ depends on \mathbf{p} . However, it does not depend on the order of the p_j .



We can consider directly the ordered vector

$$\mathbf{p} \in \nabla_\infty = \{(p_1, p_2, \dots), p_1 \geq p_2 \geq \dots > 0, \sum p_i = 1\}.$$

For instance, p_3 = the frequency of the third most frequent DNA type in Nature.

Bayesian nonparametrics: we need a prior for the parameter \mathbf{p} .

Bayesian nonparametrics: we need a prior for the parameter \mathbf{p} .
Two parameter Poisson Dirichlet distribution.

Bayesian nonparametrics: we need a prior for the parameter \mathbf{p} .

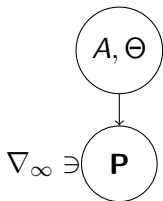
Two parameter Poisson Dirichlet distribution.

Parameters:

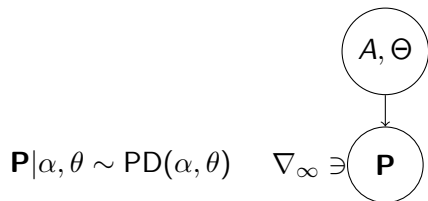
$$0 < \alpha < 1, \theta > -\alpha$$

The model (first part)

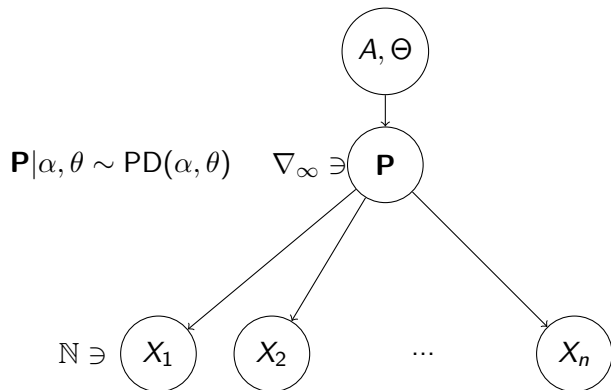
The model (first part)



The model (first part)

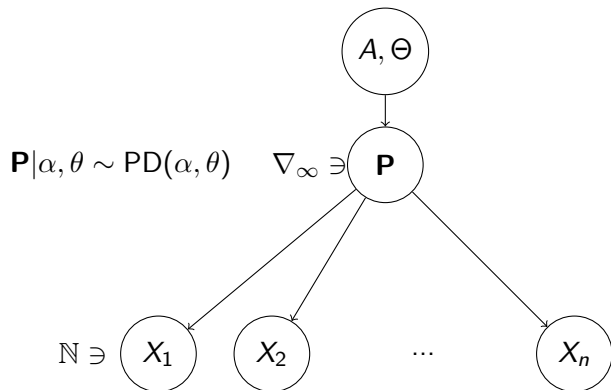


The model (first part)



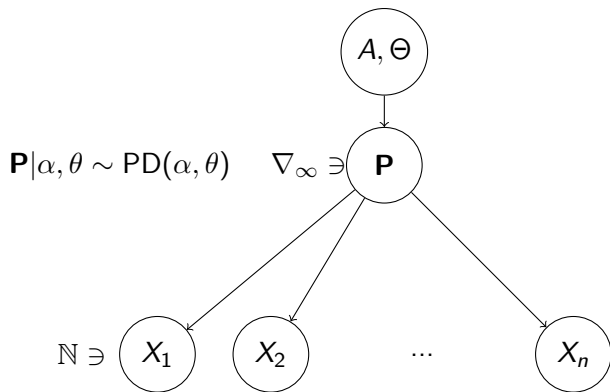
$X_i = j \rightarrow$ the i -th observation has the j th most common type in Nature.

The model (first part)

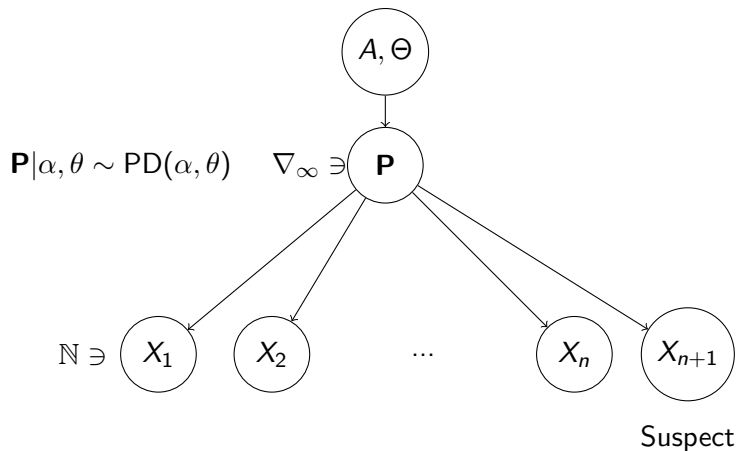


$X_i = j \rightarrow$ the i -th observation has the j th most common type in Nature.

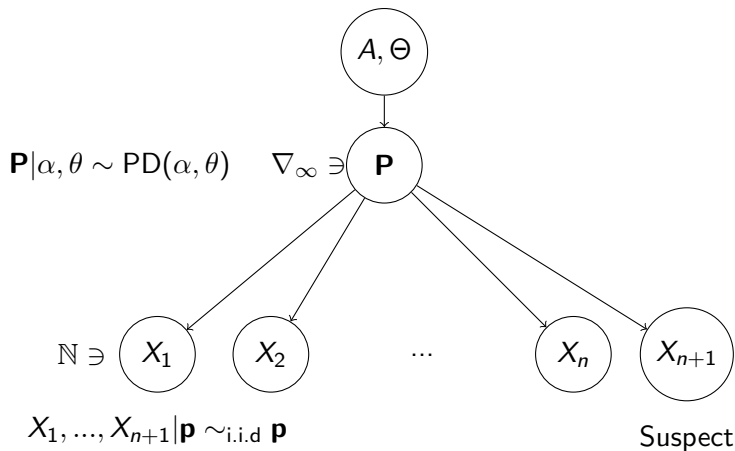
The model (first part)



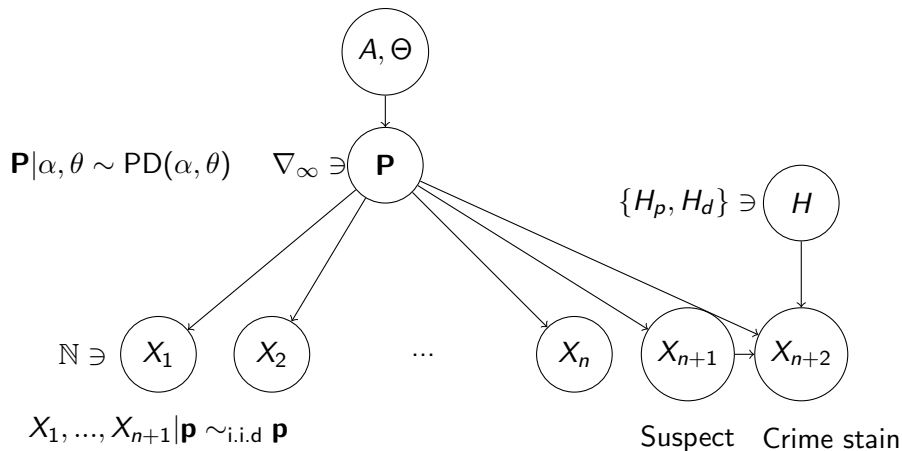
The model (first part)



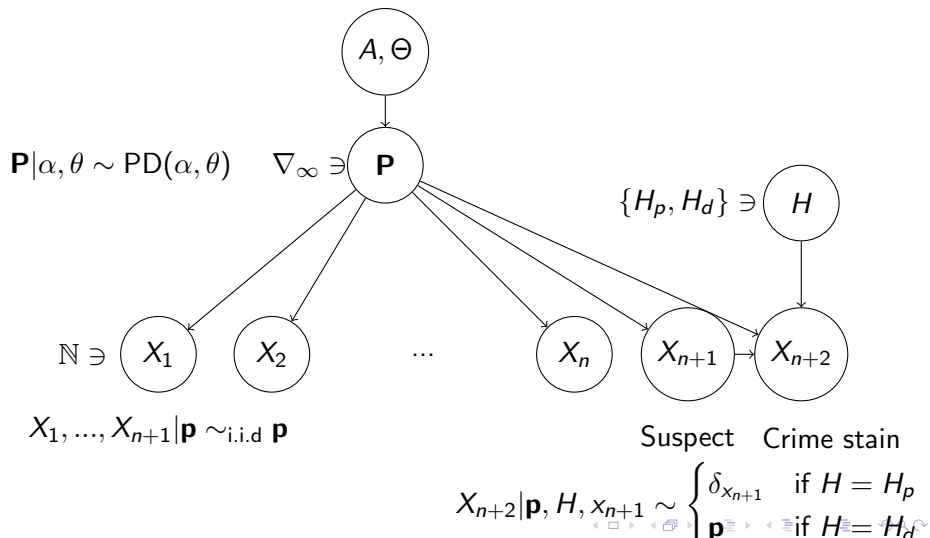
The model (first part)



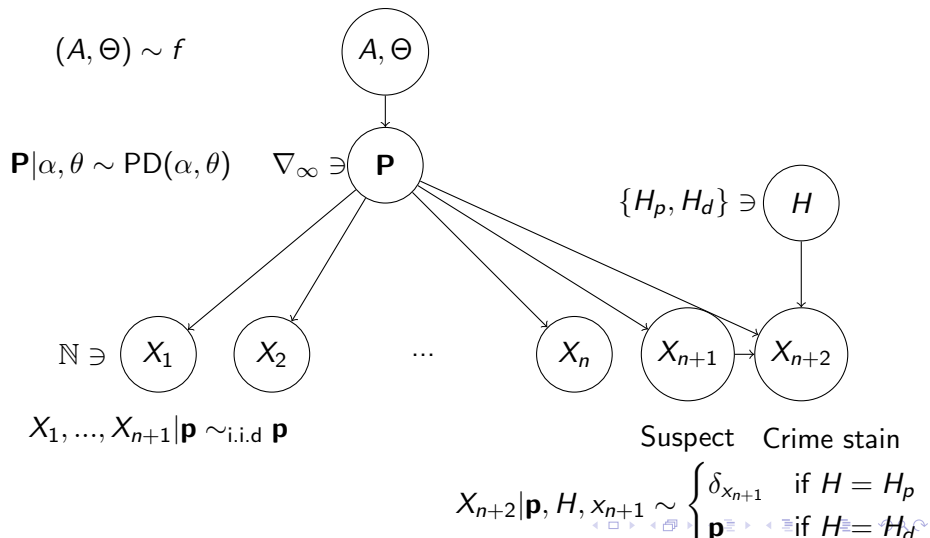
The model (first part)



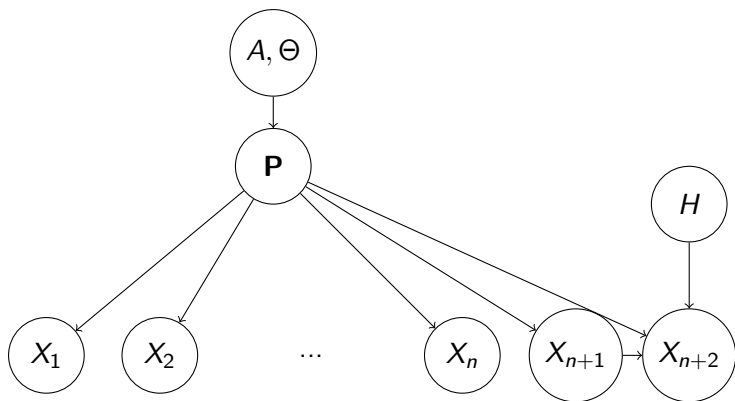
The model (first part)



The model (first part)



The model (first part)



Random partitions

Some notation:

Given $X_1, \dots, X_n \in \mathbb{N}$, random variables, $\Pi_{[n]}(X_1, X_2, \dots, X_n)$ is the random partition defined by the equivalence classes of $i \sim j$ iff $X_i = X_j$.

Random partitions

Some notation:

Given $X_1, \dots, X_n \in \mathbb{N}$, random variables, $\Pi_{[n]}(X_1, X_2, \dots, X_n)$ is the random partition defined by the equivalence classes of $i \sim j$ iff $X_i = X_j$.

$$X_1, \dots, X_n \longrightarrow \Pi_{[n]} = \pi_{[n]}^{\text{Db}}$$

$$X_1, \dots, X_n, X_{n+1} \longrightarrow \Pi_{[n+1]} = \pi_{[n+1]}^{\text{Db}+}$$

$$X_1, \dots, X_n, X_{n+1}, X_{n+2} \longrightarrow \Pi_{[n+2]} = \pi_{[n+2]}^{\text{Db}++}$$

Random partitions

Some notation:

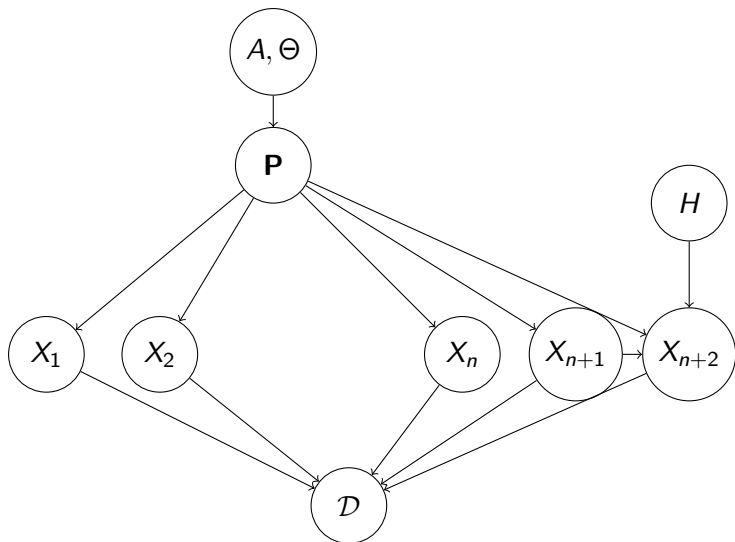
Given $X_1, \dots, X_n \in \mathbb{N}$, random variables, $\Pi_{[n]}(X_1, X_2, \dots, X_n)$ is the random partition defined by the equivalence classes of $i \sim j$ iff $X_i = X_j$.

$$\begin{array}{lll} X_1, \dots, X_n & \longrightarrow & \Pi_{[n]} = \pi_{[n]}^{\text{Db}} \\ X_1, \dots, X_n, X_{n+1} & \longrightarrow & \Pi_{[n+1]} = \pi_{[n+1]}^{\text{Db}+} \\ X_1, \dots, X_n, X_{n+1}, X_{n+2} & \longrightarrow & \Pi_{[n+2]} = \pi_{[n+2]}^{\text{Db}++} \end{array}$$

X_1, \dots, X_n are not observed, but generates the same partition as the original database.

Data can be defined as $\mathcal{D} = \Pi_{[n+2]}$.

The complete model



Pitman sampling formula

Pitman sampling formula

$$\mathbf{P} \sim \text{PD}(\alpha, \theta)$$

$$X_1, X_2, \dots, X_n | \mathbf{P} = \mathbf{p} \sim_{\text{i.i.d}} \mathbf{p}$$

Pitman sampling formula

$$\mathbf{P} \sim \text{PD}(\alpha, \theta)$$

$$X_1, X_2, \dots, X_n | \mathbf{P} = \mathbf{p} \sim_{\text{i.i.d}} \mathbf{p}$$

then $\Pi_{[n]} = \Pi_{[n]}(X_1, \dots, X_n)$ has the following distribution:

Pitman sampling formula

$$\mathbf{P} \sim \text{PD}(\alpha, \theta)$$

$$X_1, X_2, \dots, X_n | \mathbf{P} = \mathbf{p} \sim_{\text{i.i.d}} \mathbf{p}$$

then $\Pi_{[n]} = \Pi_{[n]}(X_1, \dots, X_n)$ has the following distribution:

$$\Pr(\Pi_{[n]} = \pi_{[n]} | \alpha, \theta) = \mathbb{P}_{\alpha, \theta}^n(\pi_{[n]}) = \frac{[\theta + \alpha]_{k-1; \alpha}}{[\theta + 1]_{n-1; 1}} \prod_{i=1}^k [1 - \alpha]_{n_i - 1; 1},$$

Pitman sampling formula

$$\mathbf{P} \sim \text{PD}(\alpha, \theta)$$

$$X_1, X_2, \dots, X_n | \mathbf{P} = \mathbf{p} \sim_{\text{i.i.d}} \mathbf{p}$$

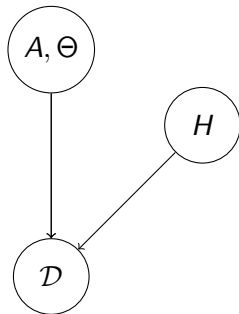
then $\Pi_{[n]} = \Pi_{[n]}(X_1, \dots, X_n)$ has the following distribution:

$$\Pr(\Pi_{[n]} = \pi_{[n]} | \alpha, \theta) = \mathbb{P}_{\alpha, \theta}^n(\pi_{[n]}) = \frac{[\theta + \alpha]_{k-1; \alpha}}{[\theta + 1]_{n-1; 1}} \prod_{i=1}^k [1 - \alpha]_{n_i-1; 1},$$

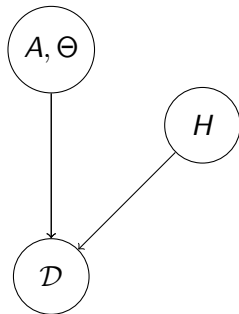
In our model

$$\Pr(D | \alpha, \theta, h) = \Pr(\Pi_{[n+2]} = \pi_{[n+2]}^{Db++} | \alpha, \theta, h) = \begin{cases} \mathbb{P}_{\alpha, \theta}^{n+2}(\pi_{[n+2]}^{Db++}) & \text{if } h = H_d \\ \mathbb{P}_{\alpha, \theta}^{n+1}(\pi_{[n+1]}^{Db+}) & \text{if } h = H_p \end{cases}$$

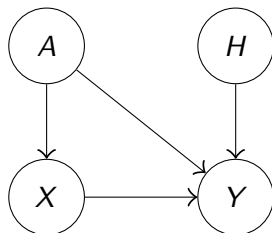
The model, simplified



The model, simplified



$$\mathcal{D} = \Pi_{[n+2]}.$$

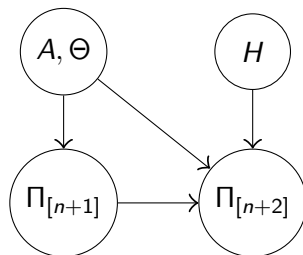


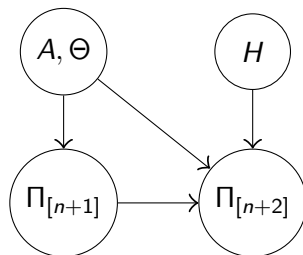
Lemma

Given four random variables A , H , X and Y , as above, the likelihood function for h , given $X = x$ and $Y = y$, satisfies

$$\text{lik}(h \mid x, y) \propto \mathbb{E}(p(y \mid x, A, h) \mid X = x).$$

Lemma





$$\text{lik}(h \mid \pi_{[n+1]}, \pi_{[n+2]}) \propto \mathbb{E}(p(\pi_{[n+2]} \mid \pi_{[n+1]}, A, \Theta, h) \mid \Pi_{[n+1]} = \pi_{[n+1]}).$$

$$\text{LR} = \frac{p(\pi_{[n+2]}|H_p)}{p(\pi_{[n+2]}|H_d)} = \frac{p(\pi_{[n+1]}, \pi_{[n+2]}|H_p)}{p(\pi_{[n+1]}, \pi_{[n+2]}|H_d)} = \frac{\text{lik}(H_p|\pi_{[n+1]}, \pi_{[n+2]})}{\text{lik}(H_d|\pi_{[n+1]}, \pi_{[n+2]})}$$

$$\text{LR} = \frac{p(\pi_{[n+2]}|H_p)}{p(\pi_{[n+2]}|H_d)} = \frac{p(\pi_{[n+1]}, \pi_{[n+2]}|H_p)}{p(\pi_{[n+1]}, \pi_{[n+2]}|H_d)} = \frac{\text{lik}(H_p|\pi_{[n+1]}, \pi_{[n+2]})}{\text{lik}(H_d|\pi_{[n+1]}, \pi_{[n+2]})}$$

Lemma allows to write

$$\text{LR} = \frac{\mathbb{E}\left(\overbrace{p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_p)}^1 \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}{\underbrace{\mathbb{E}\left(p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_d) \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}_{\frac{1-A}{n+1+\Theta}}}$$

$$\text{LR} = \frac{p(\pi_{[n+2]}|H_p)}{p(\pi_{[n+2]}|H_d)} = \frac{p(\pi_{[n+1]}, \pi_{[n+2]}|H_p)}{p(\pi_{[n+1]}, \pi_{[n+2]}|H_d)} = \frac{\text{lik}(H_p|\pi_{[n+1]}, \pi_{[n+2]})}{\text{lik}(H_d|\pi_{[n+1]}, \pi_{[n+2]})}$$

Lemma allows to write

$$\begin{aligned} \text{LR} &= \frac{\mathbb{E}\left(\overbrace{p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_p)}^1 \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}{\mathbb{E}\left(\overbrace{p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_d)}^{\frac{1-A}{n+1+\Theta}} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)} \\ &= \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}. \end{aligned}$$

$$\text{LR} = \frac{p(\pi_{[n+2]}|H_p)}{p(\pi_{[n+2]}|H_d)} = \frac{p(\pi_{[n+1]}, \pi_{[n+2]}|H_p)}{p(\pi_{[n+1]}, \pi_{[n+2]}|H_d)} = \frac{\text{lik}(H_p|\pi_{[n+1]}, \pi_{[n+2]})}{\text{lik}(H_d|\pi_{[n+1]}, \pi_{[n+2]})}$$

Lemma allows to write

$$\begin{aligned} \text{LR} &= \frac{\mathbb{E}\left(\overbrace{p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_p)}^1 \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}{\mathbb{E}\left(\overbrace{p(\pi_{[n+2]} | \pi_{[n+1]}, A, \Theta, H_d)}^{\frac{1-A}{n+1+\Theta}} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)} \\ &= \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}. \end{aligned}$$

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

By defining the random variable $\Phi = n \frac{1-A}{n+1+\Theta}$ we can write the LR as

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

By defining the random variable $\Phi = n \frac{1-A}{n+1+\Theta}$ we can write the LR as

$$\text{LR} = \frac{n}{\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]})}.$$

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

By defining the random variable $\Phi = n \frac{1-A}{n+1+\Theta}$ we can write the LR as

$$\text{LR} = \frac{n}{\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]})}.$$

We are interested in the distribution of $\Phi, \Theta \mid \Pi_{[n+1]}$

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

By defining the random variable $\Phi = n \frac{1-A}{n+1+\Theta}$ we can write the LR as

$$\text{LR} = \frac{n}{\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]})}.$$

We are interested in the distribution of $\Phi, \Theta \mid \Pi_{[n+1]}$

$$p(\phi, \theta \mid \pi_{[n+1]}) \propto p(\pi_{[n+1]} \mid \phi, \theta) f(\phi, \theta)$$

$$\text{LR} = \frac{1}{\mathbb{E}\left(\frac{1-A}{n+1+\Theta} \mid \Pi_{[n+1]} = \pi_{[n+1]}\right)}$$

By defining the random variable $\Phi = n \frac{1-A}{n+1+\Theta}$ we can write the LR as

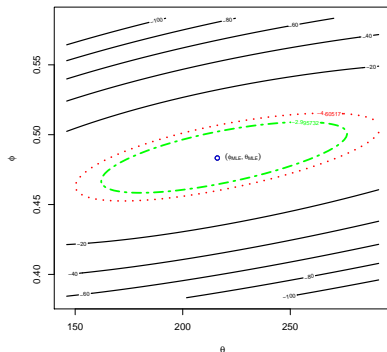
$$\text{LR} = \frac{n}{\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]})}.$$

We are interested in the distribution of $\Phi, \Theta \mid \Pi_{[n+1]}$

$$p(\phi, \theta \mid \pi_{[n+1]}) \propto p(\pi_{[n+1]} \mid \phi, \theta) f(\phi, \theta)$$

Log likelihood with ϕ and θ

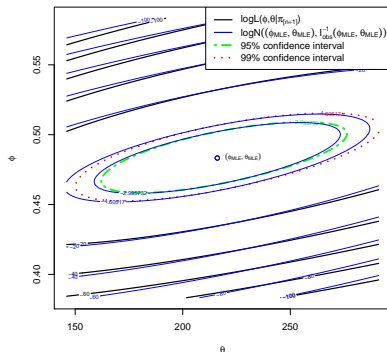
$$\log_{10} p(\pi_{[n+1]} \mid \phi, \theta)$$



Dutch Y-STR database, 7 loci, N=18,925

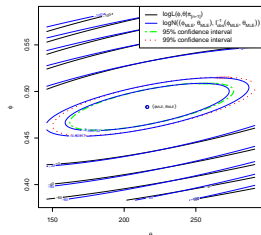
Log likelihood with ϕ and θ

$$\log_{10} p(\pi_{[n+1]} | \phi, \theta)$$



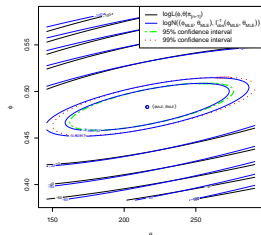
Dutch Y-STR database, 7 loci, N=18,925

Log likelihood as a function of ϕ and θ



$$p(\pi_{[n+1]} | \phi, \theta) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1})$$

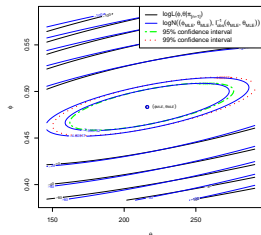
Log likelihood as a function of ϕ and θ



$$p(\pi_{[n+1]} | \phi, \theta) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1})$$

$$p(\phi, \theta | \pi_{[n+1]}) \propto p(\pi_{[n+1]} | \phi, \theta) f(\phi, \theta)$$

Log likelihood as a function of ϕ and θ

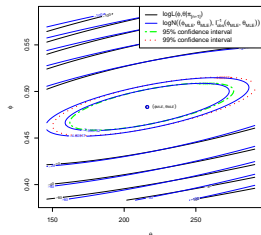


$$p(\pi_{[n+1]} \mid \phi, \theta) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1})$$

$$p(\phi, \theta \mid \pi_{[n+1]}) \propto p(\pi_{[n+1]} \mid \phi, \theta) f(\phi, \theta)$$

If the prior is smooth around the MLE then

Log likelihood as a function of ϕ and θ



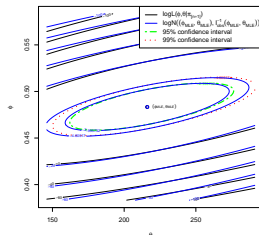
$$p(\pi_{[n+1]} \mid \phi, \theta) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1})$$

$$p(\phi, \theta \mid \pi_{[n+1]}) \propto p(\pi_{[n+1]} \mid \phi, \theta) f(\phi, \theta)$$

If the prior is smooth around the MLE then

$$p(\phi, \theta \mid \pi_{[n+1]}) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1}).$$

Log likelihood as a function of ϕ and θ



$$p(\pi_{[n+1]} \mid \phi, \theta) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1})$$

$$p(\phi, \theta \mid \pi_{[n+1]}) \propto p(\pi_{[n+1]} \mid \phi, \theta) f(\phi, \theta)$$

If the prior is smooth around the MLE then

$$p(\phi, \theta \mid \pi_{[n+1]}) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1}).$$

$$p(\phi, \theta \mid \pi_{[n+1]}) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1}).$$

$$p(\phi, \theta \mid \pi_{[n+1]}) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1}).$$

It follows that $\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]}) \approx \phi_{MLE}$.

$$p(\phi, \theta \mid \pi_{[n+1]}) \approx N((\phi_{MLE}, \theta_{MLE}), I_{MLE}^{-1}).$$

It follows that $\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]}) \approx \phi_{MLE}$.

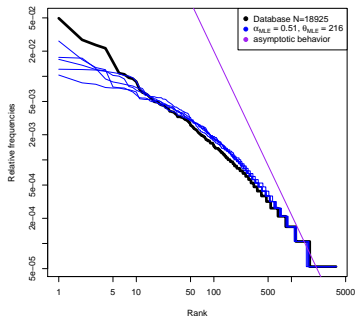
$$\text{LR} = \frac{n}{\mathbb{E}(\Phi \mid \Pi_{[n+1]} = \pi_{[n+1]})} \approx \frac{n + 1 + \theta_{MLE}}{1 - \alpha_{MLE}}$$

Sorted relative frequencies: how good is our prior?

Comparison between the spectrum from a big database, and simulations from $PD(\alpha, \theta)$ using MLE estimators of the parameters.

Sorted relative frequencies: how good is our prior?

Comparison between the spectrum from a big database, and simulations from $PD(\alpha, \theta)$ using MLE estimators of the parameters.



Thick black line: ranked relative frequencies in the database.

Thin black lines: simulations from the $PD(\alpha_{MLE}, \theta_{MLE})$.

Dotted line: asymptotics.

The LR when \mathbf{p} is known

Imagine we know \mathbf{p} .

The LR when \mathbf{p} is known

Imagine we know \mathbf{p} .

$$\text{LR}_{|\mathbf{p}} = \frac{p(\pi_{[n+2]}^{\text{Db}++} | H_p, \mathbf{p})}{p(\pi_{[n+2]}^{\text{Db}++} | H_d, \mathbf{p})} =$$

The LR when \mathbf{p} is known

Imagine we know \mathbf{p} .

$$\text{LR}_{|\mathbf{p}} = \frac{p(\pi_{[n+2]}^{\text{Db}++} | H_p, \mathbf{p})}{p(\pi_{[n+2]}^{\text{Db}++} | H_d, \mathbf{p})} = \text{Applying Lemma} =$$

The LR when \mathbf{p} is known

Imagine we know \mathbf{p} .

$$\text{LR}_{|\mathbf{p}} = \frac{p(\pi_{[n+2]}^{\text{Db}++} | H_p, \mathbf{p})}{p(\pi_{[n+2]}^{\text{Db}++} | H_d, \mathbf{p})} = \text{Applying Lemma} = \frac{1}{\mathbb{E}(p_{x_{n+1}} | \pi_{[n+1]}^{\text{Db}+}, \mathbf{p})}.$$

The LR when \mathbf{p} is known

Imagine we know \mathbf{p} .

$$\text{LR}_{|\mathbf{p}} = \frac{p(\pi_{[n+2]}^{\text{Db}++} | H_p, \mathbf{p})}{p(\pi_{[n+2]}^{\text{Db}++} | H_d, \mathbf{p})} = \text{Applying Lemma} = \frac{1}{\mathbb{E}(p_{x_{n+1}} | \pi_{[n+1]}^{\text{Db}+}, \mathbf{p})}.$$

How is this compared to the one we get with our method when \mathbf{p} is unknown?

Test Dutch database (N=2085, 7 loci)

Database of 2085 Y-STR profiles from Dutch men.

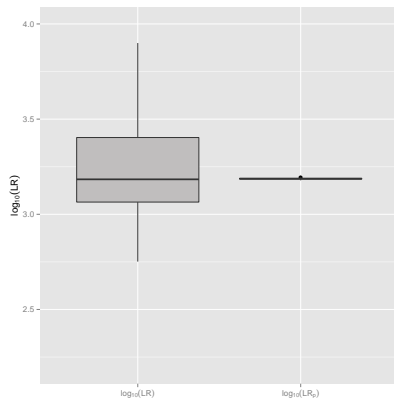
Test Dutch database (N=2085, 7 loci)

Database of 2085 Y-STR profiles from Dutch men.

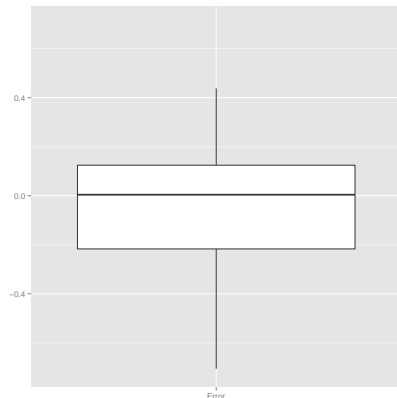
Test: Compare the distribution of $\log_{10}(LR_{|p})$ and $\log_{10} LR$ obtained by 100 samples of size 100 from this population.

Results

Compare the distribution of $\log_{10}(\text{LR}_{|p})$ and $\log_{10} \text{LR}$ obtained by 100 samples of size 100 from this population



(a) Comparison



(b) Error